



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

The Best of Both Worlds

Citation for published version:

Morris, JG 2016, The Best of Both Worlds: Linear Functional Programming without Compromise. in *ICFP 2016 Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming*. ACM SIGPLAN Notices, no. 9, vol. 51, ACM, Nara, Japan, pp. 448-461, 21st ACM SIGPLAN International Conference on Functional Programming, Nara, Japan, 18/09/16. <https://doi.org/10.1145/2951913.2951925>

Digital Object Identifier (DOI):

[10.1145/2951913.2951925](https://doi.org/10.1145/2951913.2951925)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

ICFP 2016 Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



The Best of Both Worlds

Linear Functional Programming without Compromise

J. Garrett Morris

The University of Edinburgh, UK
Garrett.Morris@ed.ac.uk

Abstract

We present a linear functional calculus with both the safety guarantees expressible with linear types and the rich language of combinators and composition provided by functional programming. Unlike previous combinations of linear typing and functional programming, we compromise neither the linear side (for example, our linear values are first-class citizens of the language) nor the functional side (for example, we do not require duplicate definitions of compositions for linear and unrestricted functions). To do so, we must generalize abstraction and application to encompass both linear and unrestricted functions. We capture the typing of the generalized constructs with a novel use of qualified types. Our system maintains the metatheoretic properties of the theory of qualified types, including principal types and decidable type inference. Finally, we give a formal basis for our claims of expressiveness, by showing that evaluation respects linearity, and that our language is a conservative extension of existing functional calculi.

Categories and Subject Descriptors D.3.2 [*Language Classifications*]: Applicative (functional) languages; D.3.3 [*Language Constructs and Features*]: Polymorphism

Keywords linear types; substructural types; qualified types

1. Introduction

Integers have a pleasing consistency: values do not become more or less integers over the course of a computation. The same is not true for file handles: we can no longer expect to read from or write to a file handle after it has been closed. Traditional functional type systems, like the logic they resemble, are good for integers (i.e., unchanging propositions), but less good for file handles (i.e., temporary ones). If our type systems are to help in the latter case, we need ones with a different logical character.

One approach is suggested by Girard’s linear logic [11], which requires that each hypothesis be used exactly once in the course of a proof. Intuitively, linear propositions are finite resources, which can neither be duplicated nor discarded, rather than arbitrary truth values, available as often, or as rarely, as needed. Linear type systems adopt the same approach to variables: each bound variable must be used exactly once in the body of its binder. Such type systems

have been used to reason about resource usage and concurrency. For example, they have been used to assure safe manipulation of state [1, 13], regulate access to shared resources [4, 7], and type interacting concurrent processes [5, 34, 36]. Each of these examples uses the restrictions on reuse and discard introduced by linearity to assure safety invariants. Simultaneously, several general purpose linear functional languages have been proposed, including those of Wadler [35], Mazurak et al. [24], and Tov and Pucella [33]. However, attempts to adapt functional programming idioms and abstractions to these calculi are thwarted by the interplay of linear and intuitionistic types. This paper proposes a novel combination of linear and qualified types that provides the safety of linear types without losing the expressiveness of functional programming.

We identify three requirements for the integration of linear types and functional programming. As an example, consider the K combinator, defined as $\lambda x. \lambda y. x$. We begin with its arguments: argument x is used once, and so can take on values of any type. Argument y is discarded, and so can only take on values of unrestricted type. This illustrates the first requirement: we must distinguish between quantification over arbitrary type and quantification over unrestricted types. Next, consider the application $K V$, giving a new function $\lambda y. V$. Whether we can reuse this function depends on V . If V is a Boolean or integer value, for example, there is no danger in reusing $\lambda y. V$. On the other hand, if V is a file handle or capability, then reuse of the function would imply reuse of V , and should be prevented. This illustrates the second requirement: we must distinguish between linear and unrestricted functions, which distinction is determined by the environment captured by each function. In particular, there is no single static characterization of the linearity of the subterms of K that accounts for its application to both linear and unrestricted values. Finally, consider the composition function \circ , defined as $\lambda f. \lambda g. \lambda x. f(g x)$. We know that we must be able to apply f and g to things, and that both are used linearly. However, these constraints are satisfied by both linear and unrestricted functions. So, the final requirement is that we must generalize the typing of application to range over the possible types of function.

These requirements have been addressed in previous work, although the interactions between them have not. Quantification over unrestricted values can be expressed in a linear system using either kinds (and subkinding) [24] or type classes [9]; both approaches extend naturally to account for pairs and sums. Many existing systems use subtyping to account for application, either implicitly [9, 24] or explicitly [10, 33]. Tov and Pucella [33] introduce a notion of relative arrow qualifiers, sufficient to express the typing of K , but at the cost of significant complexity in their type system. The interplay of these disparate mechanisms has not been fully explored. For example, none the existing systems can express the desired typing of composition, nor have they been shown to support complete type inference.

We propose a new, uniform approach to integrating linear types and functional programming, based on the theory of qualified types [17]. Rather than invent new type system features, we present a language design based on a novel combination of qualified and linear typing, both existing, well-studied type systems. To demonstrate our design, we return to the K combinator, to which we give the (qualified) type $(\text{Un } u, t \geq f) \Rightarrow t \rightarrow u \xrightarrow{f} t$. First, we observed that the second argument (here typed by u) must be unrestricted; this is captured by the predicate $\text{Un } u$. Second, we observed that the result of $K V, \lambda y.V$, may be either an linear or unrestricted function, depending on the linearity of V . We capture this through the use of two predicates, one that identifies functions and another that specifies relative linearity. The predicate $\text{Fun } f$ is satisfied only when f is a function type; we write $t \xrightarrow{f} u$ to denote the type $f t u$ under the predicate $\text{Fun } f$. Here, we use it to range over the possible types for $\lambda y.V$; we would make a similar use of the Fun predicate to express the argument types of the composition operator. The predicate $t \geq f$ is satisfied when t supports more structural rules (i.e., duplication and discard) than f . Thus, in the typing of the K combinator, if t is linear, then f must be linear; alternatively, if t is unrestricted, then f can be either linear or unrestricted.

Formally, we capture our approach in the design of a core linear calculus, which we call a Qualified Linear Language (Quill). Quill is a linear variant of Jones’s calculus OML, extended with Haskell-like first-class polymorphism [19], and with entailment rules for the Un , Fun , and \geq predicates. We preserve the metatheoretic properties of OML, particularly principal types and decidable type inference, without requiring the programmer to provide type or linearity annotations. We show that our system is a conservative extension of (non-linear) OML; concretely, this means that we can view our approach as giving linear refinements of existing functional languages and idioms, rather than replacing them entirely. Finally, we give a natural (big-step) semantics for Quill and show that evaluation respects linearity.

In summary, this paper contributes:

- The design and motivation of Quill, including examples of Quill’s application to prototypical uses of linear types (dyadic session types) and higher-order functional programming (monads) (§3).
- A formal account of the Quill type system and its relationship to OML, including a sound and complete type inference algorithm (§4).
- A linearity-aware semantics of Quill, and a proof that values of linear type are neither duplicated nor discarded during evaluation (§5).
- A discussion of further extensions of Quill, including its applicability to other substructural type systems, such as affine or relevant typing (§6).

We begin with an introduction to linear type systems and their uses (§2), and conclude by discussing related (§7) and future (§8) work.

2. Substructural Type Systems

Before describing the details of our language, we give examples of several applications of substructural type systems and several general-purpose substructural calculi.

2.1 Applications of Substructural Typing

Linear type systems restrict the use of weakening (i.e., discarding variables) and contraction (i.e., reusing variables), allowing us to reason about state and resource usage in programs. For example, excluding weakening could prevent memory or resource leaks, by

requiring that each input to a computation be consumed during its evaluation. Excluding contraction could be used in describing component layouts in circuits, where a limited number of each computational unit are available. Linear type systems combine these, providing exact control over resource usage. This section describes two uses of linear types: session types and referentially-transparent in-place update. These demonstrate two different uses of linear types: session types evolve over the course of a computation, capturing changes in underlying state, while mutable values must be used linearly to preserve referential transparency. In each case, we make two points: first, the need for the restrictions introduced by linearity, and second, the difficulties introduced in integrating the use of linear and unrestricted types.

Mutable arrays. We begin by considering in-place update. Suppose that we want to be able to read from and update arrays in a referentially transparent way. We might expect each update to produce a new copy of the array; otherwise, updates would be visible through other references to the original array. For large arrays, this copying will be extremely costly, both in time and space. The copying could be avoided if we could ensure that the use of arrays was single-threaded. That is, so long as no “old” copies of arrays are ever used, updates can be performed in place. Chen and Hudak [6] consider the connection between single-threaded usage, potentially enforced monadically, and linearity. They introduce an affine type system (they allow discarding but not duplication), and show that updating linearly-typed values can safely be performed in place. They also show how the operations on a linear data type can be interpreted to give a monad in an intuitionistic calculus, while preserving the safety of in-place update. However, this approach relies on hiding the linearly typed values, making them second-class citizens of the non-linear calculus. That is, while the monadic computations describe transformations of an underlying array, they give no direct access to the array itself. Consequently, while their approach applies to linear arrays of unrestricted element types, it could not apply to linear element types (such as other arrays), because it gives direct access to the array elements.

Session types. Next we consider session types, an instance of behavioral typing. Communication protocols frequently specify not just what data can be sent, but in what order. For example, the Simple Mail Transfer Protocol specifies not just a list of commands (identifying senders, recipients, message bodies, and so forth), but also a particular ordering to messages (the sender’s address must precede the recipients’ addresses, which must precede the message body). Session types, originally proposed by Honda [16], provide a mechanism for capturing such expectations in the types of communication channels. The critical aspect of his type system is that types evolve over the course of a computation to reflect the communication that has already taken place. For example, if channel c has session type $\text{Int} ! \text{Int} ? \text{End}$, we expect to send an integer along c , then receive an integer from c . After we have sent an integer, the type of c must change to $\text{Int} ? \text{End}$, reflecting the remaining expected behavior. We can implement session types in a functional setting by giving channels linear types, and reflecting the evolution of types in the type signatures of the communication primitives:

$$\text{send} :: t \otimes (t ! s) \rightarrow s \quad \text{receive} :: (t ? s) \rightarrow t \otimes s$$

Continuing the example above, we see that the result of $\text{send } (4, c)$ will be of type $\text{Int} ? \text{End}$, as we hoped. The linearity of these channels is crucial to assuring the type correctness of communication: reusing channel c would allow us to send arbitrarily many integers, not just one. There are approaches to encoding session types in existing functional languages, such as that of Pucella and Tov [29], but they result in channels being second class values. For example,

sending or receiving channels requires different primitives from those for sending or receiving other values.

2.2 General-Purpose Linear Calculi

Wadler [35] gives a λ -calculus based on Girard’s logic of unity, a refinement of linear logic. In his approach, the types (ranged over by τ, v) are precisely the propositions of linear logic, including pairs $(\tau \otimes v)$, functions $(\tau \multimap v)$, and the exponential modality $!(\tau)$. His type system tracks two kinds of assumptions, linear $(x : \langle \tau \rangle)$ and intuitionistic $(x : [\tau])$; only the latter are subject to contraction and weakening. Wadler does not include polymorphism in his calculus; nevertheless, we can see that his treatment of intuitionistic types would preclude attempts toward generality. He gives explicit term constructors to introduce and eliminate the exponential modality, and these constructs surround any use of unrestricted types. If M is of type $\tau \multimap v$, and N is of type τ , then we can construct the application MN of type v ; on the other hand, if M is of type $!(\tau \multimap v)$, then we must explicitly eliminate the $!$ constructor at each use of M , as $\text{let } !f = M \text{ in } fN$. Returning to our introductory example, we have two families of types (and corresponding terms) for the K combinator, $!(\tau \multimap !v \multimap !\tau)$ if the first argument is intuitionistic, and $!(\tau \multimap !v \multimap \tau)$ otherwise.

Mazurak et al. [24] present a streamlined, polymorphic linear λ -calculus. Their calculus, called F° , extends the Girard-Reynolds polymorphic λ -calculus with linearity, and introduces a kind system which distinguishes between linear (kind \circ) and unrestricted (kind \star) types. They then define the kinds of types such as pairs in terms of the kinds of their components: $\tau \otimes v$ is of kind \star if both τ and v are of kind \star , and must be of kind \circ otherwise. Finally, they introduce a subkinding relation, allowing a type of kind \star to be used any place a type of kind \circ is expected. This reflects the observation that an unrestricted value can be used any number of times, including once. While their approach seamlessly encompasses many uses of unrestricted types, it does not extend to functions. F° distinguishes between linear functions $\lambda^\circ x.M$, of type $\tau \multimap v$, which may capture arbitrary variables in their environment, and unrestricted functions $\lambda^\star x.M$, of type $\tau \multimap v$, which can only capture unrestricted values. Consequently, F° still has four distinct types for the K combinator

$$\begin{array}{ll} \forall(t : \circ). \forall(u : \star). t \multimap u \multimap t & \forall(t : \circ). \forall(u : \star). t \multimap u \multimap t \\ \forall(t : \star). \forall(u : \star). t \multimap u \multimap t & \forall(t : \star). \forall(u : \star). t \multimap u \multimap t \end{array}$$

each with distinct inhabitants. The problem is endemic to the use of higher-order functions; for example, their system has numerous distinct application and composition functions.

Tov and Pucella [33] present Alms, an affine calculus with a kind system similar to F° but with additional flexibility in the treatment of functions. Their treatment of functions includes not just affine $(\overset{\circ}{\multimap})$ and unrestricted $(\overset{\star}{\multimap})$ functions, but also functions with relative qualifiers. For example, Alms has a single most-general type for the K combinator, written

$$\forall(t : A). \forall(u : U). t \overset{u}{\multimap} u \overset{t}{\multimap} t.$$

The arrow $\overset{t}{\multimap}$ must be more restricted than the instantiation of t . If t is instantiated to an affine type, then $\overset{t}{\multimap}$ must be $\overset{\circ}{\multimap}$; otherwise, it can be $\overset{\star}{\multimap}$. They include subtyping explicitly; for example, $\text{Int} \overset{u}{\multimap} \text{Int}$ is a subtype of $\text{Int} \overset{\circ}{\multimap} \text{Int}$. Despite the (not insignificant) complexity of their system, it is still not clear that it fully supports the expressiveness of traditional functional programming languages. For example, their system has distinct composition op-

erators with the types:

$$\begin{array}{l} \forall(t : A)(u : A)(v : A). (t \overset{\circ}{\multimap} u) \overset{u}{\multimap} (u \overset{\circ}{\multimap} v) \overset{\circ}{\multimap} (t \overset{\circ}{\multimap} v) \\ \forall(t : A)(u : A)(v : A). (t \overset{\star}{\multimap} u) \overset{u}{\multimap} (u \overset{\star}{\multimap} v) \overset{\star}{\multimap} (t \overset{\star}{\multimap} v) \end{array}$$

These types are not related by the subtyping relation, as subtyping is contravariant in function arguments.

3. Programming in Quill

This section gives an intuitive overview of our calculus Quill and its primary features. We begin by describing the use of overloading to capture the non-linear use of assumptions. We then consider the particular problems arising from having both linear and unrestricted functions, the overloading of application and abstraction, and introduce the corresponding predicates on types. Finally, we consider two examples of programming in Quill: a simple presentation of dyadic session types, demonstrating the use of linearity, and a Haskell-like presentation of monads, demonstrating the interaction between linearity and higher-order functional programming. For the purposes of this section, we use a Haskell-like syntax for Quill, in which we distinguish linear functions $(\tau \overset{\circ}{\multimap} v)$ from unrestricted functions $(\tau \overset{\star}{\multimap} v)$. We give a formal account of Quill’s syntax and semantics in the following sections.

3.1 Contraction and Weakening with Class

Our goal is a functional language in which values of some (but not all) types must be treated linearly. The central problem is the integration of unrestricted types, and functions on unrestricted (but otherwise generic) types, with an otherwise linear type system. We describe one solution, based on the theory of qualified types.

We begin by distinguishing linear from unrestricted types. We consider a type to be unrestricted if values of that type can be duplicated and discarded. That is, a type τ is unrestricted if we can exhibit values of type $\tau \overset{\star}{\multimap} 1$ and $\tau \overset{\star}{\multimap} \tau \otimes \tau$. (This approach roughly parallels Filinski’s interpretation of intuitionistic types by commutative comonoids in the model of a linear calculus [8].) For example, consider a type for Booleans with the standard branching construct and constants. We can demonstrate that Booleans are unrestricted by giving the terms $\lambda b. \text{if } b \text{ then } () \text{ else } ()$ to discard a Boolean, and $\lambda b. \text{if } b \text{ then } (True, True) \text{ else } (False, False)$ to copy one. This leaves the problem of how to write code generic over such types; for instance, we would like the function $\lambda x. (x, x + 1)$ to be applicable to arguments of any unrestricted numeric type.

Our approach is inspired by the use of type classes in Haskell. Type classes were introduced to solve similar problems, such as how to write functions generic over types that have an equality operator, or that can be converted to and from text. For our purposes, we can imagine introducing a type class Un , which identifies unrestricted types:

```
class Un t where
  drop :: t -> 1
  dup  :: t -> t ⊗ t
```

The methods of Un provide the defining behavior of an unrestricted type. We could then imagine using these methods to implement terms such as the one above, for which we could write $\lambda x. \text{let } (x, x') = \text{dup } x \text{ in } (x, x' + 1)$. In inferring a type for this term, we would observe that its argument type has to support numeric operations (and so be a member of the Num class), and has to support dup (and so be a member of the Un class). We would conclude that it should have type $(\text{Num } t, \text{Un } t) \Rightarrow t \overset{\star}{\multimap} t \otimes t$.

One advantage of this view of unrestricted types is that it extends naturally to products, sums, and recursive types. For example, a pair of values (V, W) can safely be copied only when both

V and W could individually be copied. We can capture this in an instance of the `Un` class:

```
instance (Un t, Un u) => Un (t ⊗ u) where
  drop (x, y) = let () = drop x in drop y
  dup (x, y) = ((x', y'), (x'', y'')) where
    (x', x'') = dup x
    (y', y'') = dup y
```

The relationship between the linearity of t and u and the linearity of $t \otimes u$ arises organically from the typing of the `drop` and `dup` methods. The argument for sums is parallel, with the same results.

Of course, we do not intend programmers to use the `drop` and `dup` methods directly, and we imagine that instances of `Un` would be inferred automatically from type declarations. Instead, Quill allows variables to be used freely, and infers `Un` predicates as if any duplication or discarding of variables had been done explicitly. Thus, $\lambda x.(x, x + 1)$ is a well-typed Quill term with the type $(\text{Num } t, \text{Un } t) \Rightarrow t \dot{\rightarrow} t \otimes t$, as above.

3.2 The Problem of the Copyable Closure

We have an appealing view of how to distinguish unrestricted from linear types, and how to account for the linearity of products and sums. Unfortunately, this view does not extend to provide a uniform treatment of functions. Consider the curried pair constructor $\lambda x.\lambda y.(x, y)$. We know that the linearity of the resulting pair depends only on the linearity of its components. But what about the intermediate result? Suppose that we apply this function to some value V giving the term $\lambda y.(V, y)$. Whether we can copy this term depends upon the captured value V ; intuitively, we can say it depends on the function's closure. However, this is not reflected in the function type. (While the type of V does appear in the result type, so does the type of y , but the linearity of the function type is solely a consequence of V 's type.) We are thus forced to introduce distinct types for linear and unrestricted functions. This section discusses the resulting language design questions: how to handle application and abstraction in a language with multiple function types, and how to relate the type of a function to the type of its captured environment.

Application

We begin with application, the simpler of the two problems. Consider the uncurried application function. In intuitionistic calculi, this is $\lambda(f, x).f x$, of type $(t \rightarrow u, t) \rightarrow u$. In the linear settings, things are not so simple: we must decide whether the argument f and the function being defined are linear or unrestricted functions. These choices are independent, giving four incomparable types:

$$\begin{array}{ll} (t \rightarrow u) \otimes t \rightarrow u & (t \rightarrow u) \otimes t \dot{\rightarrow} u \\ (t \dot{\rightarrow} u) \otimes t \rightarrow u & (t \dot{\rightarrow} u) \otimes t \dot{\rightarrow} u \end{array}$$

We can resolve this repetition by observing that (built-in) application is implicitly overloaded: we would like to write $f x$ whether f is a linear or unrestricted function. We make this overloading explicit in the types. We introduce a new predicate, `Fun` f , which holds when f is a function type; intuitively, we can think of this as corresponding to a class whose sole method is application, and whose only members are \rightarrow and $\dot{\rightarrow}$. We can then type application with reference to this class, rather than in terms of either of the concrete function types. This reduces the number of application functions from four to two: we have $\text{Fun } f \Rightarrow f t u \otimes t \rightarrow u$ and $\text{Fun } f \Rightarrow f t u \otimes t \dot{\rightarrow} u$. We introduce syntactic sugar to make the `Fun` predicate easier to read. We will write $t \xrightarrow{f} u$ to indicate the type $f t u$ constrained by `Fun` f , and further write $t \rightarrow u$ to indicate $t \xrightarrow{f} u$ for some fresh type variable f . Using this sugar, we arrive at the most general type for the application function, $(t \rightarrow u) \otimes t \rightarrow u$.

In the previous section, we motivated the typing of contraction and weakening using the methods of an `Un` class, even though we intend their use to be implicit in practice. In the case of the `Fun` predicate, the class method intuition is less helpful. Defining primitive application as a class method is difficult (how would it be used, except by application?), and we will rely on the `Fun` predicate holding only for the built-in function types. This reinforces the expressiveness of qualified types, even beyond their traditional application to overloaded class methods.

Abstraction

We have accounted for the uncurried application function. Now consider its curried equivalent, expressed in an intuitionistic setting as $\lambda f.\lambda x.f x$ of type $(t \rightarrow u) \rightarrow t \rightarrow u$. The problem here is similar to the problem with the K combinator or the curried pair constructor. Suppose that we apply this function to some value V , giving $\lambda x.V x$: whether this function needs to be linear depends on the linearity of V . We thus have six incomparable types for the curried application function:

$$\begin{array}{lll} (t \rightarrow u) \rightarrow t \rightarrow u & (t \dot{\rightarrow} u) \rightarrow t \rightarrow u & (t \dot{\rightarrow} u) \rightarrow t \dot{\rightarrow} u \\ (t \rightarrow u) \dot{\rightarrow} t \rightarrow u & (t \dot{\rightarrow} u) \dot{\rightarrow} t \rightarrow u & (t \dot{\rightarrow} u) \dot{\rightarrow} t \dot{\rightarrow} u \end{array}$$

Our approach to overloading application allows us to give names to individual function arrows in a type. Unfortunately, even this is not sufficient to account for the types of the application function; it only allows us to reduce the six types above to two:

$$(t \xrightarrow{f} u) \rightarrow t \xrightarrow{f} u \quad (t \dot{\rightarrow} u) \rightarrow t \rightarrow u$$

However, this observation suggests our actual solution. Consider a more general type, subsuming the two above (but admitting one erroneous case): $(t \xrightarrow{f} u) \rightarrow t \xrightarrow{g} u$. The first case above is where f and g are the same type, and the second is where f is less restricted (i.e., admits more structural rules) than g . The erroneous case is where f is more restricted (i.e., admits fewer structural rules) than g . We introduce a new predicate, $\tau \geq v$, which holds when τ admits more structural rules than v . We can now give the principal type of the application operator: $f \geq g \Rightarrow (t \xrightarrow{f} u) \rightarrow t \xrightarrow{g} u$.

Our examples have focused on function types. However, the \geq relation is not limited to functions; for example, consider the possible types of the curried pair constructor $\lambda x.\lambda y.(x, y)$:

$$\text{Un } t \Rightarrow t \rightarrow u \rightarrow t \otimes u \quad t \rightarrow u \dot{\rightarrow} t \otimes u$$

As for the application operator, we see that the linearity of the final arrow is restricted by the types appearing before it in the type signature. Unlike in that case, however, the earlier type in question is not a function type. We can give the curried pair constructor the principal type $t \geq f \Rightarrow t \rightarrow u \xrightarrow{f} t \otimes u$.

3.3 Quill in Action

One of the pleasing aspects of this work has been the simplicity of our motivating examples: the K combinator and application functions are very short, but reveal the unique benefits of Quill. We conclude this section by turning to several larger examples. First, we consider a simple embedding of dyadic session types, a typical application of linear typing. Doing so demonstrates that we have not made our system too permissive. Second, we consider a presentation of Haskell's monad class and several of its instances. This shows that Quill supports the full generality of intuitionistic functional programming abstractions, and demonstrates the additional information captured by a linear type system.

For these examples, we will assume various language features present in Haskell, such as new type definitions, multi-parameter type classes with functional dependencies, and `do` notation for

monads. We believe these are representative of realistic settings for linear functional programming. However, these ideas are not fundamental to our approach, and our formalization in the following sections will consider a core calculus that does not assume such language features or syntactic sugar.

Dyadic Session Types

Session types, introduced by Honda [16], provide a typing discipline for communication protocols among asynchronous processes. There is a significant body of work exploring the combination of session-typed and functional programming. Much of this work has focused on defining new linear calculi, combining functional and concurrent programming [10, 23, 34]. These calculi frequently include details specific to session typing in their type systems, and so seem a poor fit for general purpose programming languages. Pucella and Tov [29] give an encoding of session types in Haskell, wrapping an underlying untyped use of channels. They express the session typing discipline using the existing features of the Haskell class system. However, they treat channels as second-class values, capturing the session types of channels in a parameterized monad [2] rather than in the types of the channels themselves. One consequence of this is that sending and receiving channels, while possible, requires primitive operations (with particularly involved types) distinct from those for sending and receiving values. We will show that Quill allows us to have the best of both worlds: because Quill is linear, we can have first-class channels, and because Quill fits into the existing work on qualified types we can encode the session typing discipline without having to extend our core type system.

Honda gives five constructors for session types ζ , interpreted as follows:

$\tau ! \zeta$	Send a value of type τ , then continue as ζ
$\tau ? \zeta$	Receive a value of type τ , then continue as ζ
$\zeta \boxplus \zeta'$	Choose between behaviors ζ and ζ'
$\zeta \boxtimes \zeta'$	Offer a choice of behaviors ζ and ζ'
End	No communication

(Our syntax for the choice constructors differs from Honda's to avoid conflict with the notation for the linear logic connectives.) Lindley and Morris [23] observed that, in a linear functional setting, the choice types can be encoded in terms of \oplus and the input and output types, and so we omit them from our example. We introduce types for the remaining session types—these types are empty, as we will use them as tags rather than to type channels directly.

```
data t !: s
data t ? : s
data End
```

Honda observed that communicating processes had dual expectations for their shared channels: if one process expects to send a value of type τ , the other process should expect to receive a value of type τ . Following Pucella and Tov [29], we can capture this using a type class with functional dependencies [20]:

```
class Dual t u | t → u, u → t
instance Dual s s' ⇒ Dual (t !: s) (t ? : s')
instance Dual s s' ⇒ Dual (t ? : s) (t !: s')
instance Dual End End
```

We now turn to channels and their primitive operators.

```
data Ch s
instance Un (Ch End)
```

Unlike other approaches to encoding session types in functional languages, we treat End channels as unrestricted, avoiding the need for explicit close operations. Previous work on linearity has discussed the encapsulation of unrestricted types in linear ones, either

via existential types [24, 33] or via a module system [29]. Alternatively, one might prefer to take the notion of linear channels as primitive. Either approach is possible in Quill; as we are primarily concerned with the use of linear types, we omit further discussion of them here. (But see the extended version of this paper [25] for the details of the packaging approach.) The primitive operations on session-typed channels are as follows:

```
fork    :: Dual s s' ⇒ (Ch s → M ()) → M (Ch s')
send    :: t ≥ f ⇒ t → Ch (t !: s)  $\xrightarrow{f}$  M (Ch s)
receive :: Ch (t ? : s) → M (t ⊗ Ch s)
```

We adopt the fork construct of Lindley and Morris both for its simplicity and because it assures deadlock freedom. The Dual predicate assures that the session types s and s' are well-formed and dual. Gay and Vasconcelos [10] give two typings for the send function, depending on the linearity of its first argument:

$$\begin{aligned} t &\xrightarrow{\bullet} (t ! s) \xrightarrow{\bullet} s && \text{if } t \text{ is unrestricted} \\ t &\xrightarrow{\bullet} (t ! s) \xrightarrow{\circ} s && \text{otherwise} \end{aligned}$$

This fits precisely the pattern captured by the \geq predicate in Quill. Finally, as the communication primitives are side-effecting, we assume the results are embedded in some suitable monad M . (This is not an entirely innocuous choice; we will return to monads in a linear setting for our next example.)

We present a simple example using session-typed channels. We begin with a process that performs an arithmetic operation:

```
multiplier c =
  do (x, c) ← receive c
  (y, c) ← receive c
  send (x * y) c
  return ()
```

The multiplier function defines a process that expects to read two numbers on channel c , and then sends their product back along the same channel. The inferred type for multiplier is $\text{Num } t \Rightarrow \text{Ch } (t ? : (t ? : (t !: \text{End}))) \rightarrow M ()$. Note that, despite our reuse of the name c , each call to a communication primitive returns a new copy of the channel, which is used linearly. Next, we define a process to communicate with multiplier. To illustrate the use of channels as first-class values, we define it in a round-about way. First, we define a process that provides only one of the two expected values:

```
sixSender c =
  do (d, c) ← receive c
  send 6 d
  send d c
  return ()
```

This function defines a process that begins by receiving a channel d along c ; it then sends 6 along the received channel before returning the received channel along c . Thus, its type is

$$\text{Num } t \Rightarrow \text{Ch } (\text{Ch } (t !: s) ? : \text{Ch } s !: \text{End}) \rightarrow M ()$$

Finally, we can define the main process, which uses the preceding processes to compute 42:

```
answer = do d ← fork sixSender
           c ← fork multiplier
           d ← send c d
           (c, d) ← receive d
           c ← send 7 c
           (x, c) ← receive c
           return x
```

This example demonstrates the advantages of Quill for linear programming. Unlike encoding-based approaches, we have simple types and uniform treatment of channels and other data. Unlike other concurrency-focused approaches, we have not built any aspects of session typing into our language or its type system.

Monads

In the previous example, we assumed that we could express our communication primitives monadically, to account for their side effects. As they are fundamentally reliant on higher-order functions, it is worth examining the interaction between linearity and the monadic combinators. For a simple example, consider the desugaring of answer, which begins

```
fork sixSender >>= \d → fork multiplier >>= \c → M
```

where M denotes the remainder of answer, and both c and d are free in M . As d is of linear type, we see that $\lambda c.M$ must be a linear function. Does this mean that the result of $\gg=$ must also be linear? How does this play out for other monads, like the Maybe monad?

Of course, we could transport standard intuitionistic definitions of monads directly into Quill, treating all functions as unrestricted. Doing so would allow us to use monads for unrestricted values without any new complexity. However, doing so would also rule out interesting cases, such as those with channels in the previous example. Here we take the opposite perspective, attempting to generalize standard notions of monads to include the linear cases. We will consider two canonical examples, failure and state.

First, we consider failure. We assume we have some type Maybe t with constructors Just and Nothing; observe that Maybe t is unrestricted precisely when t is unrestricted. To demonstrate that Maybe is a monad, we give implementations of the return and ($\gg=$) operators, as follows:

```
return = \x → Just x
(>>=) = \m → \f → case m of
  Nothing → Nothing
  Just x → f x
```

The typing of return is uninteresting. On the other hand, consider the use of f in the body of ($\gg=$): if m is Nothing, then f is discarded, whereas if m is Just x , then f is used once. So, we see that f must be unrestricted, and so we have the types:

```
return :: t → Maybe t
(>>=) :: t ≥ f ⇒ Maybe t → (t → Maybe u)  $\xrightarrow{f}$  Maybe u
```

The requirement that f be unrestricted captures that the remainder of the computation may not occur, an important characteristic of the failure monad. For example, this means that the monad !M! in the session types example cannot include exceptions. This should align with our expectations: if a process fails, it cannot fulfill its outstanding session-typed obligations.

Next, we consider the state monad. A state monad for state values of type S is typically implemented in Haskell by the type $S \rightarrow (t, S)$. This introduces additional choice in the linear case: should we consider values of type $S \xrightarrow{\circ} t \otimes S$ or of type $S \xrightarrow{\bullet} t \otimes S$? What constraints would this choice impose on the use of the monad? We can clarify these questions by considering the definition of return and ($\gg=$). (Relying on our generalization of abstraction and application, we consider these implementations in parallel with the choice of the state monad itself.)

```
return = \x → \s → (x, s)
(>>=) = \m → \f → \s → let (x, s') = m s in f x s
```

We make two observations about the state monad. First, x is captured in $\backslash s \rightarrow (x, s)$; therefore, a state computation can only be as unrestricted as its result values. (This is true of the failure monad as well, but is reflected in the inherent linearity of Maybe types.) Second, note that the function f is used linearly in the body of ($\gg=$), so its type need not be unrestricted (unlike for the failure monad). These observations are reflected in the types of return and ($\gg=$). We begin by introducing an alias for the state monad type:

```
type State k s t = s  $\xrightarrow{k}$  (t  $\otimes$  s)
```

Term variable	$x, y \in \text{Var}$	Type variables	$t, u \in \text{TVar}$
Multienvironments	H	Environments	Γ, Δ
Type constructors	$T^\kappa \in \mathcal{T}^\kappa$ where $\{\oplus, \bullet, \circ\} \subseteq \mathcal{T}^{* \rightarrow * \rightarrow *}$		
Kinds	$\kappa ::= * \mid \kappa \rightarrow \kappa$		
Types	$\tau^\kappa ::= t \mid T^\kappa \mid \tau^{\kappa' \rightarrow \kappa} \tau^{\kappa'}$		
Predicates	$\pi ::= \text{Un } \tau \mid \text{Fun } \tau \mid \tau \geq v$		
Qualified types	$\rho ::= \tau^* \mid \pi \Rightarrow \rho$		
Type schemes	$\sigma ::= \rho \mid \forall t. \sigma$		
Expressions	$M, N ::= x \mid KM \mid \lambda x. M \mid MN \mid \text{in}_1 M \mid \text{in}_2 N$ $\mid \text{case } M \text{ of } \{ \text{in}_1 x \mapsto N; \text{in}_2 y \mapsto N' \}$ $\mid \text{let } x = M \text{ in } N \mid \text{let } Kx = M \text{ in } N$		

Figure 1: Quill types and terms.

We can then type return and ($\gg=$) by

```
return :: t ≥ State k s t ⇒ t → State k s t
(>>=) :: (State k s t ≥ g, f ≥ State k s u) ⇒
  State k s t → (t  $\xrightarrow{f}$  State k s u)  $\xrightarrow{g}$  State k s u
```

The predicate $\text{State } k s t \geq g$ reflects that the term $\backslash f \rightarrow \dots$ has captured m of type $\text{State } k s t$.

Finally, we generalize these examples. The problem is the type of the second argument to ($\gg=$): to be useful in the linear context, we must sometimes include the restricted function type, but to incorporate the full range of monads we must sometimes limit it to unrestricted functions. We encompass both cases using a multi-parameter type class for monads:

```
class Monad f m | m → f where
  return :: t ≥ m t ⇒ t → m t
  (>>=) :: (m t ≥ g, f ≥ m u) ⇒
    m t → (t  $\xrightarrow{f}$  m u)  $\xrightarrow{g}$  m u
```

The definitions above give instances of our new Monad class:

```
instance Monad (→) Maybe
instance Monad k (State k s)
```

and that the example of dyadic session types will type in monads m such that $\text{Monad } (\xrightarrow{\circ}) m$ is provable.

We should emphasize that, because f is functionally dependent on m , our reformulation of the Monad class does not introduce any new polymorphism, or new potential for ambiguity. Rather, it makes explicit (at the type level) existing differences in the composition of monadic computations.

4. Substructural Qualified Types

We have considered some of the challenges of using linear calculi in practice, given an intuitive description of how we addresses these challenges using qualified types, and demonstrated how our solution might be realized in a Haskell-like practical programming language. In this section, we give a formal account of our approach to substructural qualified types. We begin by giving an overview of a core Quill calculus and its type system (§4.1). We then give a syntax-directed variant on the type system (§4.2), preparatory to giving an Algorithm \mathcal{M} [22] style type inference algorithm (§4.3). Finally, we relate Quill typing to typing for a non-substructural core calculus (§4.4), making concrete our claims that Quill encompasses existing functional programming practice.

4.1 Quill Terms and Typing

The syntax of Quill types and terms is shown in Figure 1. Quill types are stratified according to a simple kind system; we write

$\boxed{P \mid H \vdash M : \sigma}$			
(VAR) $\frac{}{P \mid \lambda x : \sigma \vdash x : \sigma}$	(CTR) $\frac{P \mid H, H', H' \vdash M : \sigma \quad P \vdash H' \text{ un}}{P \mid H, H' \vdash M : \sigma}$	(WKN) $\frac{P \mid H \vdash M : \sigma \quad P \vdash H' \text{ un}}{P \mid H, H' \vdash M : \sigma}$	
$\frac{P \mid H, x : \tau \vdash M : v}{P \Rightarrow \text{Fun } \phi \quad P \vdash H \geq \phi}$	(\rightarrow I) $\frac{}{P \mid H \vdash \lambda x. M : \phi \tau v}$	$\frac{P \mid H \vdash M : \phi \tau v}{P \mid H' \vdash N : \tau \quad P \Rightarrow \text{Fun } \phi}$	(\rightarrow E) $\frac{}{P \mid H, H' \vdash M N : v}$
(LET) $\frac{P \mid H \vdash M : \sigma \quad P \mid H', x : \sigma \vdash N : \tau}{P \mid H, H' \vdash \text{let } x = M \text{ in } N : \tau}$			
(PLI) $\frac{P \mid H \vdash M : \tau_i}{P \mid H \vdash \text{in}_i M \vdash \tau_1 \oplus \tau_2}$		(PE) $\frac{P \mid H \vdash M : \tau_1 \oplus \tau_2 \quad P \mid H'_1, x : \tau_1 \vdash N : v \quad P \mid H'_2, x : \tau_2 \vdash N' : v}{P \mid H, H'_x \vdash \text{case } M \text{ of } \{ \text{in}_1 x \mapsto N; \text{in}_2 x \mapsto N' \} : v}$	
(MAKE) $\frac{K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \xrightarrow{\bullet} \tau \quad P \Rightarrow [\vec{v}/\vec{u}]Q \quad P \mid H \vdash M : [\vec{v}/\vec{u}]\tau' \quad \vec{u} \notin \text{ftv}(P, H)}{P \mid H \vdash K M : \tau}$		(BREAK) $\frac{K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \xrightarrow{\bullet} \tau \quad P \mid H \vdash M : \tau \quad P, [\vec{v}/\vec{t}]Q \mid H', x : [\vec{v}/\vec{t}]\tau' \vdash N : v' \quad \vec{u} \notin \text{ftv}(P, Q, H, H', v')}{P \mid H, H' \vdash \text{let } K x = M \text{ in } N : v'}$	
(\Rightarrow I) $\frac{P, \pi \mid H \vdash M : \rho}{P \mid H \vdash M : \pi \Rightarrow \rho}$	(\Rightarrow E) $\frac{P \mid H \vdash M : \pi \Rightarrow \rho \quad P \Rightarrow \pi}{P \mid H \vdash M : \rho}$	(\forall I) $\frac{P \mid H \vdash M : \sigma \quad t \notin \text{ftv}(P, H)}{P \mid H \vdash M : \forall t. \sigma}$	(\forall E) $\frac{P \mid H \vdash M : \forall t. \sigma}{P \mid H \vdash M : [\tau/t]\sigma}$
$\boxed{P \vdash \cdot \text{ un}}$		$\boxed{P \vdash \cdot \geq \phi}$	
(UN- τ) $\frac{P \Rightarrow \text{Un } \tau}{P \vdash \tau \text{ un}}$	(UN- ρ) $\frac{P, \pi \vdash \rho \text{ un}}{P \vdash \pi \Rightarrow \rho \text{ un}}$	(\geq - τ) $\frac{P \Rightarrow \tau \geq \phi}{P \vdash \tau \geq \phi}$	(\geq - ρ) $\frac{P, \pi \vdash \rho \geq \phi}{P \vdash (\pi \Rightarrow \rho) \geq \phi}$
(UN- σ) $\frac{P, \text{Un } t \vdash \sigma \text{ un}}{P \vdash \forall t. \sigma \text{ un}}$	(UN-H) $\frac{\bigwedge_{x:\sigma \in H} P \vdash \sigma \text{ un}}{P \vdash H \text{ un}}$	(\geq - σ) $\frac{P, \text{Un } t \vdash \sigma \geq \phi}{P \vdash (\forall t. \sigma) \geq \phi}$	(\geq -H) $\frac{\bigwedge_{x:\sigma \in H} P \vdash \sigma \geq \phi}{P \vdash H \geq \phi}$

Figure 2: Typing rules.

τ , v and ϕ (without superscripts) to range over types of any kind. (Unlike Mazurak et al. [24], we use kind \star for all types, not just unrestricted ones.) We assume that $\xrightarrow{\bullet}$, $\xrightarrow{\bullet}$ and \oplus are binary type constructors, which we will write infix, corresponding to linear and unrestricted functions and additive sums. We do not include multiplicative or additive products ($\tau \otimes v$ and $\tau \& v$), as these can be encoded in terms of the other types. (These encodings depend on our overloading of abstraction for their full generality.) We allow arbitrary additional type constructors, providing other (user-defined) data types. Data types capture first-class universal and existential types, following the approach of Jones’s FCP [19]. While we have not used these features in our examples, we include them for two reasons. First, existential types are used prominently in other approaches to linear functional programming [24, 33], particularly to construct linear wrappers around unrestricted types, and so we show that Quill accommodates similar constructions. Second, existentials provide a another application of the techniques developed to account for the linearity of functions; we describe this at more length when consider extensions to our core calculus (§6). Predicates π include those necessary for our treatment of linearity (and can constrain higher-kinded types). Qualified types and type schemes are standard for overloaded Hindley-Milner calculi. We write $\forall \vec{t}. Q \Rightarrow \tau$ to abbreviate $\forall t_1 \dots \forall t_n. P_1 \Rightarrow \dots \Rightarrow P_m \Rightarrow \tau$.

Quill includes standard terms for variables, abstractions, applications, and (additive) sums. We introduce polymorphism at let bindings. First-class existential and universal types are expressed using constructors K . We assume an ambient signature mapping individual constructors K to types $\forall \vec{t}. (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \xrightarrow{\bullet} \tau$. (This type is not included in σ ; σ denotes inferable type schemes.) Con-

struction $K M$ builds a value of type τ , assuming that τ' has a suitably generic type. We insist that constructors be fully applied. Deconstruction let $K x = M$ in N eliminates such values. The introduction of first-class polymorphism through data types corresponds to common practice in Haskell, and allows us to make clear the extent of type inference.

Figure 2 gives the Quill type system. The typing judgment is $P \mid H \vdash M : \sigma$, where P is a set of predicates on the type variables in the remainder of the judgment, H is a typing environment, M is a term and σ a type scheme. Our treatment of the typing environment follows a standard approach for linear logic, but differs from some of the existing work on linear type systems. We use multisets of typing assumptions (which we call “multienvironments”); thus, the multienvironment $\{x : \sigma, x : \sigma\}$ is distinct from $\{x : \sigma\}$. We require that multienvironments be consistent, so if $\{x : \sigma, x : \sigma'\} \subseteq H$ then we must have $\sigma = \sigma'$, and we write H, H' for the consistent multiset union of H and H' . Assumptions of unrestricted type are duplicated in (CTR) and discarded in (WKN); both rules use the auxiliary judgment $P \vdash \cdot \text{ un}$, which lifts the Un predicate to typing environments. (The assumption Un t in (UN- σ) accounts for terms like the empty list, which should be treated as unrestricted until their type variables are instantiated.) This allows us to simplify the other rules: in (VAR), there must only be one binding in the environment, while rules like (\rightarrow E) can split the typing environment among their hypotheses. In particular, we avoid introducing an auxiliary judgment to split type environments while sharing assumptions of unrestricted type present in many linear type systems [24, 33]. Rules (\rightarrow I) and (\rightarrow E) implement overloading of abstraction and application. In

$\frac{P \ni \pi}{P \Rightarrow \pi}$	$\frac{\bigwedge_{\pi \in Q} P \Rightarrow \pi}{P \Rightarrow Q}$	$\frac{\tau = \overset{\circ}{\rightarrow} \vee \tau = \overset{\bullet}{\rightarrow}}{P \Rightarrow \text{Fun } \tau}$
$K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \overset{\bullet}{\rightarrow} \tau \quad P, Q, \text{Un } \vec{t} \Rightarrow \text{Un } \tau'$		
$P \Rightarrow \text{Un } \tau$		
$\frac{}{P \Rightarrow \text{Un } (\tau \overset{\bullet}{\rightarrow} v)}$	$\frac{P \Rightarrow \text{Un } \tau_1 \quad P \Rightarrow \text{Un } \tau_2}{P \Rightarrow \text{Un } (\tau_1 \oplus \tau_2)}$	
$\frac{P \Rightarrow \text{Un } \tau}{P \Rightarrow \tau \geq (v \overset{\bullet}{\rightarrow} v')}$	$\frac{}{P \Rightarrow \tau \geq (v \overset{\circ}{\rightarrow} v')}$	
$\frac{P \Rightarrow \tau \geq \phi t \quad t \text{ fresh}}{P \Rightarrow \tau \geq \phi}$	$\frac{P \Rightarrow \tau t \geq \phi \quad t \text{ fresh}}{P \Rightarrow \tau \geq \phi}$	

Figure 3: Entailment rules.

(\rightarrow E), note that we allow the function term to be of any type $\phi \tau v$, so long as it satisfies the constraint $\text{Fun } \phi$. In (\rightarrow I), we allow a term $\lambda x.M$ to have any function type, so long as that type is more restricted than its environment; the auxiliary judgment $P \vdash \cdot \geq \cdot$ lifts the \geq predicate to type environments. We will assume throughout this presentation that binders introduce fresh names. First-class polymorphism is introduced in (MAKE) and eliminated in (BREAK); our approach follows Jones [19] almost exactly, but adds a predicate context Q . We write $K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \rightarrow \tau$ to denote an instantiation of the signature for K in the ambient context. We assume that each data type has at most one constructor; more complex data types can be expressed using the other features of the type system. The remaining rules are standard for linear sums and qualified polymorphism.

Figure 3 gives a minimal definition of the predicate entailment relation $P \Rightarrow Q$. One strength of type systems based on qualified types is that the predicate system provides a natural point of extension, and our approach here is no different. Nevertheless, we specify some rules for the entailment judgment, namely the linearity of the built-in types and that (only) $\overset{\circ}{\rightarrow}$ and $\overset{\bullet}{\rightarrow}$ are in class Fun . In determining the linearity of a data type τ , we assume that the universally quantified type variables \vec{t} are unrestricted (as a term of type τ cannot have made any assumptions of \vec{t}), but cannot do so for the existentially quantified variables \vec{u} , as they may have been instantiated arbitrarily in constructing the τ value. We have intentionally given a minimal specification of \geq ; in particular, we have omitted various simplification rules which might be expected in a practical implementation, and have limited the our attention to cases of \geq with function types on their right-hand side (as those are the only such predicates introduced by our typing rules). We will return to the definition of this class when we discuss extensions to Quill (§6). Finally, the lifting cases for \geq are a notational convenience; for example, they allow us to write $\tau \geq \phi$ rather than $\tau \geq \phi t u$ for fresh t and u .

4.2 A Syntax-Directed Quill Type System

The Quill type system has a number of rules that are not syntax directed, including the structural rules and the rules introducing and eliminating polymorphism. To simplify the definition of type inference and the proofs of its correctness, we give a syntax-directed variant of the Quill type system. In doing so, we address two independent concerns. First, the rules (\forall I), (\forall E), (\Rightarrow I), and (\Rightarrow E) may be used at any point in a derivation. This problem has already been studied in the general context of qualified types. An identical

solution applies in Quill: uses of (\forall E) and (\Rightarrow E) may always be permuted to occur at occurrences of (\forall I), while uses of (\forall I) and (\Rightarrow I) may always be permuted to occur at occurrences of (LET) or at the end of the derivation. Second, the structural rules (WKN) and (CTR) may also appear at any point in a typing derivation. As in the polymorphism cases, we show that uses of these rules can be permuted to definite places in the derivation: uses of (CTR) can be permuted to appear immediately below a rule with multiple hypotheses (such as (\rightarrow E) or (\oplus E)) and uses of (WKN) can be permuted to occurrences of (\forall I).

Figure 4 gives the syntax-directed variant of the Quill system. The judgment $P \mid \Gamma \vdash^S M : \tau$, is a syntax-directed variant of $P \mid H \vdash M : \sigma$, and uses standard type environments Γ rather than multienvironments H . We write Γ, Γ' to denote partitioning the environment; the treatment of contraction and weakening is explicit in the typing rules, rather than via a partitioning relation on typing environments. The auxiliary judgments are unchanged. The syntax-directed system differs from the original type system in two ways. First, we account for polymorphism. Our approach is identical to Jones’s approach for (intuitionistic) qualified types [17]: we introduce instantiation and generalization operators, accounting for the role of predicates, and collapse the treatment of polymorphism into the instances of (VAR^S) and (LET^S).

Definition 1. We define instantiation and generalization as follows:

1. Let σ be some type scheme $\forall \vec{t}. P \Rightarrow \tau'$. We say that $Q \Rightarrow \tau$ is an instance of σ , written $(Q \Rightarrow \tau) \sqsubseteq \sigma$, if there is some \vec{v} such that $\tau = [\vec{v}/\vec{t}]\tau'$ and $Q \Rightarrow [\vec{v}/\vec{t}]P$.
2. Let Γ be a typing environment, and ρ a qualified type. We define $\text{Gen}(\Gamma, \rho)$ to be the type scheme $\forall (fiv(\rho) \setminus fiv(\Gamma)). \rho$.

We use instantiation in (VAR^S), collapsing a use of (VAR) and subsequent uses of (\forall E) and (\Rightarrow E), and generalization in (LET^S) collapsing a use of (LET) and preceding uses of (\Rightarrow I) and (\forall I). Second, we account for contraction and weakening. In (VAR^S), we allow an arbitrary environment, so long as the unused assumptions Δ are unrestricted. In ($\rightarrow \text{E}^S$), we partition the input environment into three parts: Γ is used exclusively in typing M , Γ' is used exclusively in typing N , and Δ is used in both; consequently, assumptions in Δ must be unrestricted. The remaining rules follow the same pattern.

The goal of the syntax-directed type system is a one-to-one correspondence between syntactic forms and typing rules. However, it is not the case that a typeable term has exactly one syntax-directed typing derivation. For example, while the contents of the linear environments are determined by the term structure, the contents of the unrestricted environments are not. For another example, consider the term $(\lambda x.x)y$. We can choose to type the abstraction as either $t \overset{\circ}{\rightarrow} t$ or $t \overset{\bullet}{\rightarrow} t$ (or even as $t \overset{f}{\rightarrow} t$ assuming the predicate $\text{Fun } f$). Each of these choices would make the term well-typed, and we assume that the terms obey the same laws (i.e., the choice of type introduces no observable distinction in the semantics of the term).

We now relate our original and syntax-directed type systems. We start with environments. Intuitively, the syntax-driven system introduces contraction when needed, guarded by Un constraints; however, a multienvironment H could contain multiple instances of assumptions with linear types. We introduce a notion of approximation between multienvironments H and environments Γ that holds when the only repeated assumptions in H are for unrestricted types.

Definition 2. If H is a multienvironment, Γ is an environment, and P is some context, then we say that Γ approximates H under P , written $P \vdash H \approx \Gamma$, if $x : \sigma \in H$ if and only if $x : \sigma \in \Gamma$, and if $\lambda x : \sigma, x : \sigma \sqsubseteq H$, then $P \vdash \sigma \text{ un}$.

$(\text{VAR}^S) \frac{P \vdash \Delta \text{un} \quad (P \Rightarrow \tau) \sqsubseteq \sigma}{P \mid \Delta, x : \sigma \vdash^S x : \tau}$		$(\text{LET}^S) \frac{Q \mid \Gamma, \Delta \vdash^S M : \tau \quad \sigma = \text{Gen}(\Gamma, \Delta; Q \Rightarrow \tau) \quad P \mid \Gamma', \Delta, x : \sigma \vdash^S N : v \quad P \vdash \Delta \text{un}}{P \mid \Gamma, \Gamma', \Delta \vdash^S \text{let } x = M \text{ in } N : v}$	
$(\rightarrow \text{I}^S) \frac{P \mid \Gamma, x : \tau \vdash^S M : v \quad P \Rightarrow \text{Fun } \phi \quad P \vdash \Gamma \geq \phi}{P \mid \Gamma \vdash^S \lambda x. M : \phi \tau v}$		$(\rightarrow \text{E}^S) \frac{P \mid \Gamma, \Delta \vdash^S M : \phi \tau v \quad P \mid \Gamma', \Delta \vdash^S N : \tau \quad P \Rightarrow \text{Fun } \phi \quad P \vdash \Delta \text{un}}{P \mid \Gamma, \Gamma', \Delta \vdash^S M N : v}$	
$(\oplus \text{I}_i^S) \frac{P \mid \Gamma \vdash^S M : \tau_i}{P \mid \Gamma \vdash^S \text{in}_i M : \tau_1 \oplus \tau_2}$		$(\oplus \text{E}^S) \frac{P \mid \Gamma, \Delta \vdash^S M : \tau_1 \oplus \tau_2 \quad P \vdash \Delta \text{un} \quad P \mid \Gamma', \Delta, x : \tau_1 \vdash^S N : v \quad P \mid \Gamma', \Delta, x : \tau_2 \vdash^S N' : v}{P \mid \Gamma, \Gamma', \Delta \vdash^S \text{case } M \text{ of } \{ \text{in}_1 x \mapsto N; \text{in}_2 x \mapsto N' \} : v}$	
$(\text{MAKE}) \frac{K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \xrightarrow{\bullet} \tau \quad P \Rightarrow [\vec{v}/\vec{u}]Q \quad P \mid \Gamma \vdash M : [\vec{v}/\vec{u}]\tau' \quad \vec{t} \notin \text{fv}(P, \Gamma)}{P \mid \Gamma \vdash K M : \tau}$		$(\text{BREAK}) \frac{K : (\forall \vec{t}. \exists \vec{u}. Q \Rightarrow \tau') \xrightarrow{\bullet} \tau \quad P \mid \Gamma, \Delta \vdash M : \tau \quad P \vdash \Delta \text{un} \quad P, [\vec{v}/\vec{t}]Q \mid \Gamma', \Delta, x : [\vec{v}/\vec{t}]\tau' \vdash N : v' \quad \vec{u} \notin \text{fv}(P, Q, \Gamma, \Gamma', \Delta, v')}{P \mid \Gamma, \Gamma', \Delta \vdash \text{let } K x = M \text{ in } N : v'}$	

Figure 4: Syntax-directed typing rules.

We now turn to our primary results. First, derivations in the syntax-directed system correspond to derivations in the original system.

Theorem 3 (Soundness of \vdash^S). *If $P \mid \Gamma \vdash^S M : \tau$ and $P \vdash H \approx \Gamma$, then $P \mid H \vdash M : \tau$.*

The proof is by structural induction on the derivation of $P \mid \Gamma \vdash^S M : \tau$, and relies on introducing instances of the structural and polymorphism rules.

Second, we show completeness of the syntax directed system. A derivation in the original system may end with uses of $(\Rightarrow \text{I})$ or $(\forall \text{I})$, moving predicates from the context to the type or quantifying over free type variables. In contrast, there are no such steps in a derivations in the syntax-directed system. To account for this difference, we introduce a notion of qualified type schemes, again following Jones [17].

Definition 4. A qualified type scheme $(P \mid \sigma)$ pairs a type scheme σ with a set of predicates P . Let σ be $\forall \vec{t}. Q \Rightarrow \tau$ and σ' be $\forall \vec{t}'. Q' \Rightarrow \tau'$. We say that $(P \mid \sigma)$ is an instance of $(P' \mid \sigma')$, written $(P \mid \sigma) \sqsubseteq (P' \mid \sigma')$ iff there are \vec{v} such that $\tau = [v_i/t'_i]\tau'$ and $P, Q \Rightarrow P', [v_i/t'_i]Q'$. We treat type schemes σ as abbreviations for qualified type schemes $(\emptyset \mid \sigma)$.

We can now state the completeness of the syntax-directed system.

Theorem 5 (Completeness of \vdash^S). *If $P \mid H \vdash M : \sigma$ and $P \vdash H \approx \Gamma$, then there are some Q and τ such that $Q \mid \Gamma \vdash^S M : \tau$ and $(P \mid \sigma) \sqsubseteq \text{Gen}(\Gamma, Q \Rightarrow \tau)$.*

Intuitively, this states that for any derivation in our original type system, there is a derivation of at least as general a result in the syntax-directed system. The proof is by induction on the derivation of $P \mid H \vdash M : \sigma$. The interesting cases rely on the role of generalization and instantiation in the syntax-directed type system and the safe movement of structural rules up derivation trees.

4.3 Type Inference for Quill

Having defined a suitable target type system, we can give a type inference algorithm for Quill. We have three separate concerns during type inference. First, we use a standard Hindley-Milner treatment of polymorphism. Second, we introduce Un predicates for

non-linear use of variables. We track the variables used in each expression, and so detect when variables are reused or discarded. Third, we account for first-class polymorphism. We introduce a distinction between rigid and flexible type variables; only the latter are bound in unification. These three concerns add apparent complexity to the type inference algorithm, but can be understood separately.

The inference algorithm is given in Figure 5, in the style of Algorithm \mathcal{M} [22]. The inputs include the environment Γ , expression M and expected type τ , along with the current substitution S and the rigid type variables X . The output includes the generated predicates P , the resulting substitution S' , and a set of used (term) variables Σ . We let u_i range over fresh type variables, and let U, R, S range over substitutions. We will look at illustrative cases of the algorithm in detail; the remaining cases are constructed along the same lines.

In the variable case, we are given both the variable x and its expected type τ . We unify x 's actual type, given by Γ , with its expected type τ . This illustrates the primary difference between Algorithm \mathcal{M} and Milner's Algorithm \mathcal{W} : unification is moved as close to the leaves as possible. We defer the details of unification to a separate algorithm $Mgu_X(\tau, v)$, where the type variables in X are not bound in the resulting unification procedure. The implementation of unification does not differ from previous presentations, such as Jones's unification algorithm for FCP [19]. We return any predicates in the type scheme of x , the updated substitution, and the observation that x has been used.

The application case demonstrates the sets of used variables. We check the subexpressions M and N ; to account for the overloading of functions, we only assume that M has type $u_1 u_2 \tau$, for some function type u_1 . The variables used in M are captured by Σ , and those used in N are captured by Σ' . Any variables used in both must be unrestricted, and so the predicates inferred for the application include not just the predicates inferred for each sub expression (P and P'), but also that any variables used in $\Sigma \cap \Sigma'$ must have unrestricted type. We capture this with the auxiliary function $\text{Un}(\Gamma|_{\Sigma \cap \Sigma'})$, where $\Gamma|_{\Sigma}$ denotes the restriction of Γ to variables in Σ . We give a declarative specification of $\text{Un}(-)$; an implementation that finds the simplest such P can be straightforwardly derived from the definitions of $P \vdash \cdot \text{un}$ and entailment.

The let case demonstrates the treatment of polymorphism and binders. First, we must account for the possibility that x was not used in N , and thus must be of unrestricted type. This is captured

$\mathcal{M}(S, X; \Gamma \vdash M : \tau) = P, S', \Sigma$	
$\mathcal{M}(S, X; \Gamma \vdash x : \tau) = ([\vec{u}/\vec{t}]/P), U \circ S, \{x\}$ where $(x : \forall \vec{t}. P \Rightarrow v) \in S\Gamma$ $U = Mgu_X([\vec{u}/\vec{t}]/v, \tau)$ $\mathcal{M}(S, X; \Gamma \vdash \lambda x. M : \tau) = Q, S', \Sigma \setminus x$ where $P, S'; \Sigma = \mathcal{M}(Mgu_X(\tau, u_1 u_2 u_3) \circ S, X;$ $\Gamma, x : u_2 \vdash M : u_3)$ $Q = \{\text{Fun } u_1\} \cup \text{Leq}(u_1, \Gamma _{\Sigma}) \cup$ $\text{Weaken}(x, u_2, \Sigma)$ $\mathcal{M}(S, X; \Gamma \vdash M N : \tau) = Q, R', \Sigma \cup \Sigma'$ where $P, R, \Sigma = \mathcal{M}(S, X; \Gamma \vdash M : u_1 u_2 \tau)$ $P', R', \Sigma' = \mathcal{M}(R, X; \Gamma \vdash N : u_2)$ $Q = P \cup P' \cup \{\text{Fun } u_1\} \cup$ $Un(\Gamma _{\Sigma \cap \Sigma'})$ $\mathcal{M}(S, X; \Gamma \vdash \text{in}_i M : \tau) = P, R, \Sigma$ where $P, R, \Sigma = \mathcal{M}(Mgu_X(\tau, u_1 \oplus u_2) \circ S, X; \Gamma \vdash M : u_i)$ $\mathcal{M}(S, X; \Gamma \vdash \text{case } M \text{ of } \{\text{in}_1 x \mapsto N; \text{in}_2 y \mapsto N'\} : \tau) =$ $(P_M, P_N, P_{N'}, Q), R'', \Sigma_M \cup \Sigma_N \cup \Sigma_{N'}$ where $P_M, R, \Sigma_M = \mathcal{M}(S, X; \Gamma \vdash M : u_1 \oplus u_2)$ $P_N, R', \Sigma_N = \mathcal{M}(R, X; \Gamma, x : u_1 \vdash N : \tau)$ $P_{N'}, R'', \Sigma_{N'} = \mathcal{M}(R', X; \Gamma, y : u_2 \vdash N' : \tau)$ $\Sigma' = (\Sigma_N \setminus \Sigma_{N'}) \cup (\Sigma_{N'} \setminus \Sigma_N) \cup$ $\Sigma_M \cap (\Sigma_N \cup \Sigma_{N'})$ $Q = Un(\Gamma _{\Sigma'}) \cup \text{Weaken}(x, u_1, \Sigma_N) \cup$ $\text{Weaken}(y, u_2, \Sigma_{N'})$	$\mathcal{M}(S, X; \Gamma \vdash \text{let } x = M \text{ in } N : \tau) = (P' \cup Q), R', \Sigma \cup (\Sigma' \setminus x)$ where $P, R, \Sigma = \mathcal{M}(S, X; \Gamma \vdash M : u_1)$ $\sigma = \text{GenI}(R\Gamma, R(P \Rightarrow u_1))$ $P', R', \Sigma' = \mathcal{M}(R, X; \Gamma, x : \sigma \vdash N : \tau)$ $Q = Un(\Gamma _{\Sigma \cap \Sigma'}) \cup \text{Weaken}(x, \sigma, \Sigma')$ $\mathcal{M}(S, X; \Gamma \vdash K M : \tau) = P, R \circ U, \Sigma$ where $K : \forall \vec{t}_1. (\forall \vec{t}_2. \exists \vec{t}_3. v') \xrightarrow{\bullet} v$ $U = Mgu_X(v, [\vec{u}_1/\vec{t}_1]\tau)$ $P, R, \Sigma = \mathcal{M}(U \circ S, X \cup \vec{t}_2; \Gamma \vdash M : [\vec{u}_1/\vec{t}_1, \vec{u}_3/\vec{t}_3]v')$ $\vec{t}_2 \# \text{fv}((R \circ U)\Gamma)$ $\mathcal{M}(S, X; \Gamma \vdash \text{let } K x = M \text{ in } N : \tau) = Q, R', \Sigma_M \cup (\Sigma_N \setminus x)$ where $K : \forall \vec{t}_1. (\forall \vec{t}_2. \exists \vec{t}_3. v') \xrightarrow{\bullet} v$ $P_M, R, \Sigma_M = \mathcal{M}(S, X; \Gamma \vdash M : v)$ $P_N, R', \Sigma_N = \mathcal{M}(R, X \cup \vec{t}_3; \Gamma, x : [\vec{u}_1/\vec{t}_1, \vec{u}_2/\vec{t}_2]v' \vdash N : \tau)$ $\vec{t}_3 \# \text{fv}(R'\Gamma, R'\tau)$ $Q = P_M \cup P_N \cup \text{Weaken}(x, \sigma, \Sigma_N) \cup$ $Un(\Gamma _{\Sigma_M \cap \Sigma_N})$
$\text{Leq}, \text{Un}, \text{Weaken}, \text{GenI}$	
$\text{Leq}(\phi, \Gamma) = \bigcup \{P \mid P \vdash \phi \leq \tau\}$ $Un(\Gamma) = \bigcup \{P \mid (y : \sigma) \in \Gamma, P \vdash \sigma \text{ un}\}$ $\text{Weaken}(x, \sigma, \Sigma) = \begin{cases} P & \text{if } x \notin \Sigma, P \vdash \sigma \text{ un} \\ \emptyset & \text{otherwise} \end{cases}$ $\text{GenI}(\Gamma, P \Rightarrow \tau) = \forall (\text{fv}(SP, \tau)). SP \Rightarrow \tau$ where S improves $\text{fv}(P) \setminus \text{fv}(\Gamma, \tau)$ in P	

Figure 5: Type inference algorithm \mathcal{M} . We let u_i range over fresh variables, and write $A \# B$ to require that A and B be disjoint.

by $\text{Weaken}(x, \sigma, \Sigma)$. Second, we consider generalization. Recall the term $(\lambda x. x) y$, where y has type τ . The algorithm will infer that this term has type τ under the assumption $\text{Fun } u$ for some variable u . But u appears neither in the typing environment nor in the result, so naively generalizing this expression would give the (apparently ambiguous) type scheme $(\text{Fun } u) \Rightarrow \tau$. However, this is not a real ambiguity: we have no way of observing the choice of u in the resulting expression, so we could assume it to be $\xrightarrow{\bullet}$ without decreasing the expressiveness or safety of type inference. We formalize this observation using an adaptation of Jones's notion of improvement for qualified types [18]. An improving substitution for a qualified type $P \Rightarrow \tau$ is a substitution S such that any satisfiable instance of $P \Rightarrow \tau$ is also a unambiguous satisfiable instance of $S(P \Rightarrow \tau)$. For example, $[\xrightarrow{\bullet}/f]$ is an improving substitution for $(\text{Un } f, \text{Fun } f) \Rightarrow \tau$, as the only ways to prove $\text{Fun } f$ are if f is $\xrightarrow{\bullet}$ or $\xrightarrow{\circ}$ and only the former is unrestricted. In the type $(\text{Fun } f) \Rightarrow \tau$, where f is not free in τ or the environment, we can instantiate f to either $\xrightarrow{\bullet}$ or $\xrightarrow{\circ}$ and cannot observe the choice. As this choice does not introduce ambiguity, we consider $[\xrightarrow{\bullet}/f]$ to be an improving substitution in such cases. We say that S is an improving substitution for X in P if S is the union of such improvements for each variable in X , and apply such an improving substitution before generalizing. Again, we give a declarative specification of $\text{GenI}(-, -)$, as the derivation of its implementation is entirely straightforward.

We can now relate type inference and the syntax-directed type system. First, inference constructs valid typings.

Theorem 6 (Soundness of \mathcal{M}). *If $\mathcal{M}(S, X; \Gamma \vdash M : \tau) = P, S', \Sigma$, then $S' P \mid S' (S\Gamma) \vdash M : S' (S\tau)$.*

The proof is by induction on the structure of M ; each case involves comparing the predicates generated in inference to the predicates needed for typing. In combination with Theorem 3, this gives a similar soundness result for inference with respect to the original type system. Next, we want to show that any valid typing can be found by inference.

Theorem 7 (Completeness of \mathcal{M}). *If S is a substitution and X is a set of type variables such that $P \mid S\Gamma \vdash^S M : S\tau$, and $S|_X = \text{id}$, then $\mathcal{M}(\text{id}, X; \Gamma \vdash M : \tau) = Q, S', \Sigma$ such that $(P \Rightarrow S\tau) \sqsubseteq \text{GenI}(S'\Gamma, S'Q \Rightarrow S'\tau)$.*

The proof is by induction on the typing derivation, observing in each case that the computed type generalizes the type in the derivation. Again, in combination with Theorem 5, we have a completeness result for inference with respect to the original type system. Finally, this allows us to give a constructive proof that Quill enjoys principal types.

Theorem 8 (Principal Types). *If $P_0 \mid H \vdash M : \sigma_0$ and $P_1 \mid H \vdash M : \sigma_1$ then there is some σ such that $\emptyset \mid H \vdash M : \sigma$ and $(P_0 \mid \sigma_0) \sqsubseteq \sigma, (P_1 \mid \sigma_1) \sqsubseteq \sigma$.*

The soundness of inference tells us that, if there are any typings for a term in an environment, then the inference algorithm will compute some typing for that term. The completeness of inference

tells us that the computed type will be at least as general as the original types.

4.4 Conservativity of Typing

We have claimed that Quill is as expressive as functional languages without linearity. To formalize that claim, we will show that any expression typeable in OML, Jones’s core calculus for qualified types [17], is also typeable in Quill.

OML is a Core ML-like language with qualified types. Its types and terms are pleasingly simple: the former contains functions, type variables, and qualified and quantified types, and the latter contains variables, applications, abstractions, and `let` (to introduce polymorphism). We do not give a full description of OML typing here, partly as it is so similar to Quill typing. In particular, as in Quill typing, OML has a syntax directed typing judgment $P \mid \Gamma \vdash_{\text{OML}}^S M : \tau$, where P is a collection of predicates, Γ an OML typing environment and τ an OML type.

The crux of our argument is that (by construction) the syntax-directed typing rules of Quill can each be seen as generalizations of the corresponding rules of OML. For example, rules (\rightarrow^S) and (\rightarrow^E^S) can introduce generalization over function types, a feature not present in OML. However, they need not do so; if all functions are unrestricted, (\rightarrow^S) and (\rightarrow^E^S) are elaborate restatements of the corresponding rules of OML. The remaining difference is in the treatment of variables: Quill may insist on predicates to capture their unrestricted use, where there are no corresponding predicates required by OML, but this will never cause a term to be ill-typed.

Theorem 9. *If $P \mid \Gamma \vdash_{\text{OML}}^S M : \tau$, then there is some Q such that $Q \mid \Gamma \vdash^S M : \tau$, and $Q \Rightarrow P$.*

OML also has a sound and complete type inference algorithm, and principal types. Thus, we see that if OML type inference accepts a given term, then Quill type inference will also accept the term, and in each case will compute its most general typing. We might hope to show the converse as well; however, we do not know of a non-linear core calculus that matches the exact features of Quill, including both qualified types and data type-mediated first class polymorphism.

We do not suggest that terms are given the same types in each setting: for example, the function $\lambda x. \lambda y. y$ is given the type $\forall t u. t \rightarrow u \rightarrow u$ in OML, whereas it would given the type $\forall t u. \text{Un } t \Rightarrow t \rightarrow u \rightarrow u$ in Quill. Similarly, sums must be taken as primitive in Quill (as their elimination form shares its environment), whereas they can be encoded in OML. Finally, as demonstrated in our earlier discussion of monads (§3.3), Quill may suggest more refined abstractions than are present in non-linear languages. Nevertheless, this result does show a strong connection between programming in Quill and programming in traditional functional languages, one which is not shared by other combinations of linear and functional programming.

5. Semantics

We motivated the discussion of linear type systems by considering examples like session types and mutable arrays, in which we wanted to avoid duplicating or discarding values of linear types. The Quill type system, however, only restricts the use of assumptions, and says nothing about the use of values directly. Further, Quill differs from other substructural calculi in several ways, including the use of overloading and the form of first-class polymorphism. In this section, we demonstrate that Quill assures that the use of values, not just of assumptions, is consistent with their typing. To do so, we define a natural semantics for Quill terms, annotated with the values introduced and eliminated in the course of evaluation. We can then show that any values used non-linearly

$Ix \ni j, k$ $\text{Value} \ni V, W ::= K^j V \mid \lambda^j x. M \mid \text{in}_i^j V$ $M, N ::= V \mid \dots$	
<hr/> $\frac{j \text{ fresh}}{\lambda x. M \Downarrow_{\emptyset}^{\lambda^j x. M} \lambda^j x. M} \quad \frac{M \Downarrow_E^I V \quad j \text{ fresh}}{K M \Downarrow_E^{I, K^j V} K^j V}$	
<hr/> $\frac{M \Downarrow_E^I \lambda^j x. M' \quad N \Downarrow_{E'}^{I'} V \quad [V/x] M' \Downarrow_{E''}^{I''} W}{M N \Downarrow_{E, E', E'', \lambda^j x. M'}^{I, I', I''} W}$	
<hr/> $\frac{M \Downarrow_E^I V \quad j \text{ fresh}}{\text{in}_i M \Downarrow_E^{I, \text{in}_i^j V} \text{in}_i^j V} \quad \frac{M \Downarrow_E^I V \quad [V/x] N \Downarrow_{E'}^{I'} W}{\text{let } x = M \text{ in } N \Downarrow_{E, E'}^{I, I'} W}$	
<hr/> $\frac{M \Downarrow_E^I \text{in}_i^j V \quad [V/x_i] N_i \Downarrow_{E'}^{I'} W}{\text{case } M \text{ of } \{ \text{in}_1 x_1 \mapsto N_1; \text{in}_2 x_2 \mapsto N_2 \} \Downarrow_{E, E', \text{in}_i^j V}^{I, I'} W}$	

Figure 6: Linearity-aware semantics for Quill.

have unrestricted type. Our approach is strongly inspired by that used by Mazurak et al. to prove a similar property of their F° calculus [24].

We begin by defining a notion of values for Quill. Intuitively, we might expect values to be abstractions, sums of values, or constructors applied to values. However, our intended safety property requires that we distinguish different instances of syntactically-identical values. The top of Figure 6 gives an extended syntax of Quill, in which values are tagged with indices from some index set Ix . The semantics will tag new values with fresh indices, and we will then rely on the indices to establish identity when showing safety. We identify a subset of values, LinVals , as linear:

$$\text{LinVals} = \{V \in \text{Value} \mid \text{if } \vdash V : \tau, \text{ then } \emptyset \not\Rightarrow \text{Un } \tau\}.$$

Our goal is to show that values in LinVals are neither duplicated nor discarded during evaluation. The contents of LinVals depend on the signatures of the constructors. For a simple example, suppose that we have some K with signature $(\exists u. u) \rightarrow T$. To show $P \Rightarrow \text{Un } T$, we would have to show that $P \Rightarrow \text{Un } u$ (where $u \notin \text{fv}(P)$). This is clearly impossible, so $K V \in \text{LinVals}$ for any value V . On the other hand $\lambda x. x$ is not in LinVals , as it can be given unrestricted type.

The bottom of Figure 6 gives a natural semantics for Quill. The evaluation relation $M \Downarrow_E^I V$ denotes that M evaluates to V ; the annotations I and E are multisets of values, I capturing all values introduced during the evaluation and E capturing all values eliminated during the evaluation. (We track values, rather than just indices, so that we can state the type safety theorem below.) Functions evaluate to themselves, but annotated with a fresh index. The only value introduced is the function, and no values are eliminated. The other introduction forms are similar, but must account for the evaluation of their subexpressions. We use call-by-value evaluation; as observed by Mazurak et al. [24], call-by-name and call-by-need evaluation may result in discarding linearly typed values during evaluation. The values introduced in evaluating an application are those introduced in evaluating each of its subexpressions and in evaluating the substituted result of the application. The values eliminated are those eliminated in each hypothesis and the function itself. The `let` and `case` rules are similar.

We can now state our desired safety property. Intuitively, if $M \Downarrow_E^I V$, we expect that each linear value introduced during evaluation (that is, each $W \in I \cap \text{LinVals}$) will appear either exactly

once, either in E or as a subexpression of the result V . For any expression M , we define $SExp(M)$ to be the subexpressions of M , defined in the predictable fashion, $Exp(M) = \mathcal{M} \dot{\cup} SExp(M)$, and $Val(M) = Exp(M) \cap Value$.

Theorem 10 (Type safety). *Let M be a closed term such that $\vdash M : \forall t. P \Rightarrow \tau$ and $M \Downarrow_E^I V$.*

1. $P \mid \emptyset \vdash V : \tau$.
2. Let $E' = E \cup Val(V)$, and let $D = I \setminus E'$ (the values discarded during evaluation) and $C = E' \setminus I$ (the values copied during evaluation). Then, $W \in D \cup C$ only if $W \notin LinVals$.

The proof is by induction over the structure of M . The key observation is that duplication and discarding can happen only as the result of substitution, and thus that linearity of variables (i.e., assumptions) is enough to assure that only unrestricted values are duplicated. We believe this argument can be straightforwardly generalized to small-step semantics, again following Mazurak et al. [24]

6. Extensions

We describe three extensions of Quill, showing the generality and flexibility of our approach.

Quill has a linear type system, in which both contraction (duplication) and weakening (discard) are limited to unrestricted types. Several alternative substructural logics exist: relevant logics, for example, exclude weakening but not contraction, and affine logics exclude contraction but not weakening. Some systems, such as that of Ahmed et al. [1] and Gan et al. [9] provide linear, affine, relevant, and unrestricted types simultaneously. Finally, there have been several type systems that introduce similar partitioning of assumptions to control side-effects, starting from Reynolds' work on Idealized Algol [30] and continuing with modern work on bunched implication [27] and separation logic. We have focused on the linear case in particular because various examples, such as session types, require its restrictions on both contraction and weakening. Nevertheless, we believe the Quill approach would apply equally well in these other cases. For example, the type system we have given has a single predicate, $\text{Un } \tau$, used both when assumptions are duplicated and when they are discarded. Alternatively, we could introduce distinct predicates for these cases, say $\text{Dup } \tau$ and $\text{Drop } \tau$. We could then redefine our existing predicate $\text{Un } \tau$ as the conjunction of $\text{Dup } \tau$ and $\text{Drop } \tau$. As in the systems of Ahmed et al. and Gan et al., we would require four arrow types. However, the remainder of the Quill approach would adapt seamlessly. We could extend the $\text{Fun } \tau$ predicate and the \geq relation to include the new arrow types, and the resulting system would continue to enjoy principal types and complete type inference.

The treatment of functions differs from the other primitive types (like products and sums) because the linearity of a function from τ to v cannot be determined from linearity of τ and v . A similar observation can be made of existential types. For example, suppose that we have two constructors with signatures $K_1 :: (\exists u. u) \dot{\rightarrow} T_1 t$ and $K_2 :: (\exists u. \text{Un } u \Rightarrow u) \dot{\rightarrow} T_2 t$. Assumptions of type $T_1 \tau$ will always be treated as linear, and assumptions of type $T_2 \tau$ always unrestricted, regardless of the choice of τ . We could view T_1 and T_2 as instances of a general T pattern (i.e., as the satisfying instances of a predicate $T t$) just as we view $\dot{\rightarrow}$ and $\dot{\rightarrow}$ as instances of a general \rightarrow pattern. Following the approach taken for functions, we would extend the \geq relation to include the T types, asserting that $\cdot \Rightarrow \tau \geq T_1 v$ and $\text{Un } \tau \Rightarrow \tau \geq T_2 v$. We would then introduce a generalized constructor $K :: (T t, u \geq t) \Rightarrow u \dot{\rightarrow} t$, and a generalized deconstructor $unK :: (T t, t \geq f) \Rightarrow t \rightarrow (\forall u. u \rightarrow r) \xrightarrow{f} r$. As the goal is generalizing over the predicate $\text{Un } u$, the body of the deconstructor cannot rely on its presence. This example

demonstrates the flexibility of Quill; in particular, it shows that our treatment of functions is an instance of a more general pattern, itself expressible in Quill. We suspect that this treatment of existentials in non-linear functional languages could be expressed in linear calculi. On the other hand, there are cases for which this approach would not be appropriate, such as the use of existentials to enforce linear use of unrestricted primitives. We believe that more practical experience would be required to determine how, and how often, this generalization of existentials should be applied.

We have treated the dup and drop methods as providing a helpful intuition for the use of the Un predicate, but have assumed that their explicit use or implementation is not of interest. However, there are cases in which providing non-trivial implementations of these methods could be useful. For example, many operating system resources, such as file handles, need to be explicitly freed. One could imagine capturing such resources as affine types in a language based on Quill, in which the drop method freed the underlying resource. Similarly, given suitable primitives, one could imagine using drop and dup to implement a kind of reference-counting scheme for resources, in which dup incremented the reference count and drop decremented it. This approach would generalize the various scope based mechanisms for managing such resources in languages such as C# and Java. The derived definitions of drop and dup for products and sums (§3.1) would extend to this setting as well. However, this would introduce a new concern: the placement of automatic inserted calls to the drop and dup methods. For a simple example, imagine that some variable y is in scope in the expression $\lambda x. M$, but not free in M . We must insert a call to $\text{drop } y$. Our current syntax-directed approach could be interpreted as moving calls to drop to the leaves of the typing derivation, but in this case that would delay the discard of y until the function $\lambda x. M$ is invoked, which might be undesirable.

7. Related Work

The past thirty years have seen a wealth of work on linear types and their applications. We summarize some of the work most directly related to our own.

In introducing substructural type systems (§2.2), we described several other general purpose calculi, including F° of Mazurak et al. [24], and Alms of Tov and Pucella [33]. These systems were both influential on the development of Quill. Our work differs from theirs in two regards. First, we have generalized the treatment of functions, and thus increased the expressiveness of function combinators. We believe that, especially given the importance of combinator-based idioms in functional programming, this is a significant advance in the usability of linear functional calculi. Second, our treatment relies on qualified types, rather than building notions of subkinding, subtyping, and variance into the type system itself. While this may seem to simply be trading one kind of complexity for another, we believe that qualified types are an independently useful language feature (a claim borne out by the experience of Haskell). Finally, we believe that qualified types are a natural way to express relationships among types, as demonstrated by our generalization of relative linearity to encompass existential types.

F° , Alms , and Quill all rely on identifying a collection of types as unrestricted (through kind mechanisms in the first cases and type predicates in ours). There are several other mechanisms to integrate linear and unrestricted types. Wadler [35] and Barber and Plotkin [3] give calculi based directly on the propositions and proofs of linear logic, in which each linear type τ has an intuitionistic counterpart $!\tau$. These calculi also draw a distinction between intuitionistic and linear assumptions, where only the former are subject to contraction and weakening. While these calculi have close logical connections, the manipulation of the $!$ modality adds

significant syntactic bureaucracy, and does not provide an obvious route to generalizing linear and unrestricted behavior. Walker [37] and Ahmed et al. [1] present systems of annotations for linearity (albeit without polymorphism). These approaches seem less well suited for programming with linear types, however. For example, they provide linear Booleans (of little expressive value, as the duplication and discarding operations for Booleans can be easily defined) and unlimited session-typed channels (presumably an empty type). Finally, they require all types to be annotated with linearity (or usage) annotations, which is acceptable in a core language but unsuited to languages used by humans. Clean adopts a similar annotation-based approach in its uniqueness typing system [31]. However, the aims of uniqueness typing and linearity are dual: in Clean, unique values can become non-unique (at the cost of some of their operations), while in a linear type system we must guarantee linearity (but can use unrestricted values linearly).

Gustavsson and Svenningsson [12] describe a system of usage annotations and bounded usage polymorphism; Hage et al [14] describe an alternative approach to usage inference based on effect typing and subeffects. These approaches differ from linear type systems in two ways. First, they treat usage separately from types; as we argued in the last paragraph, this produces confusing or empty types, like linear Booleans or unrestricted channels. Second, linearity is prescriptive, while usage types are descriptive. This means that usage can be approximated, where approximations of linearity would either lose safety or expressiveness. Consequently, usage analyses can be invisible to the programmer, whereas linearity must (to some degree) be programmer-visible. Nevertheless, our \leq predicate seems similar to approaches to subeffecting for usage and strictness [15], and we believe that investigating this similarity is valuable future work.

Finally, there have been numerous substructural approaches to typing for imperative and low-level languages, including region types [38], alias types [32], adoption and focus [7], and linear types for locations [26], and several generalizations of linear typing, including coeffect systems [28]. These approaches have similar goals to our work—establishing safety guarantees beyond those expressed in traditional type systems—but differ in their underlying calculi and do not share our focus on principality and type inference. Nevertheless, some of the ideas of these systems could be profitably applied in ours. For example, some adaptation of the adoption and focus mechanisms could avoid the rebinding present in cases such as our dyadic session types example. We think exploring the overlap of our system and the problems they address, such as exploring explicit memory management in a Quill-like language, will be important future work.

8. Future Work

We have presented Quill, a new linear functional language achieving both the safety guarantees made possible by linear types and the expressiveness of conventional functional programming languages. We have demonstrated several examples of linear and functional programming in Quill. We have shown that Quill has principal types and decidable type inference, that it is a conservative extension of existing functional calculi, and that reduction preserves linearity. We have also shown several simple extensions of the core Quill calculus, incorporating more flexible treatment of existentials, and other notions of substructural typing. We conclude by discussing several directions for future work.

We intend Quill to provide a foundation for practical functional programming with linear types. This can be tested in two ways. First, we intend to explore abstractions for linear programming. We hope to draw on existing mechanisms, such as adoption and focus [7] in the imperative setting and parameterized monads [2] in the functional setting, while taking advantage of Quill’s first-

class treatment of linearity to express these mechanisms within the language. Second, we hope to build larger programs in Quill, taking advantage of linearity to enhance safety properties; domains like concurrency and low-level programming seem particularly suited to such an approach.

Quill distinguishes between linear and unrestricted functions for type safety reasons, while the (high-level) semantics we give treats abstractions identically. Similar distinctions are drawn by usage type systems, for efficiency reasons. We believe that similar efficiency gains could be obtained in compiling Quill programs. As our treatment of functions is general, similar approaches could be applied to other types for efficiency reasons as well. For example, rather than requiring that arrays be treated linearly, we could overload the array operations to apply to both linear and unrestricted arrays, but to use efficient in-place operations when arrays were used linearly.

The central technical problem addressed by Quill is the multiplication of function spaces. Linear type systems are not the only context in which this can occur. Similar multiplications happen, for example, in type and effect systems or in systems that distinguish pointed and unpointed types [21]. We believe that the approach taken in Quill would generalize to such cases as well; in particular, we believe that tracking pointedness could be relevant in many of the same application domains in which linearity is relevant.

Acknowledgments. I thank James Cheney for feedback on an early draft of this paper, and James McKinna, Sam Lindley, and my other colleagues for helpful discussions. This work was funded by EPSRC grant number EP/K034413/1.

References

- [1] A. J. Ahmed, M. Fluet, and G. Morrisett. A step-indexed model of substructural state. In *Proceedings of the 10th ACM SIGPLAN International Conference on Functional Programming, ICFP 2005, Tallinn, Estonia, September 26-28, 2005*, pages 78–91, 2005.
- [2] R. Atkey. Parameterised notions of computation. *J. Funct. Program.*, 19(3-4):335–376, 2009.
- [3] A. Barber and G. Plotkin. Dual intuitionistic linear logic. Technical Report LFCS-96-347, University of Edinburgh, 1996.
- [4] J. Boyland. Checking interference with fractional permissions. In *Static Analysis, 10th International Symposium, SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings*, pages 55–72, 2003.
- [5] L. Caires and F. Pfenning. Session types as intuitionistic linear propositions. In *CONCUR*. Springer, 2010.
- [6] C. Chen and P. Hudak. Rolling your own MADT - A connection between linear types and monads. In *Conference Record of POPL’97: The 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium, Paris, France, 15-17 January 1997*, pages 54–66, 1997.
- [7] M. Fähndrich and R. DeLine. Adoption and focus: Practical linear types for imperative programming. In *Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI), Berlin, Germany, June 17-19, 2002*, pages 13–24, 2002.
- [8] A. Filinski. Linear continuations. In *Conference Record of the Nineteenth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Albuquerque, New Mexico, USA, January 19-22, 1992*, pages 27–38, 1992.
- [9] E. Gan, J. A. Tov, and G. Morrisett. Type classes for lightweight substructural types. In *Proceedings Third International Workshop on Linearity, LINEARITY 2014, Vienna, Austria, 13th July, 2014*, pages 34–48, 2014.
- [10] S. J. Gay and V. T. Vasconcelos. Linear type theory for asynchronous session types. *J. Funct. Program.*, 20(1):19–50, 2010.
- [11] J. Girard. Linear logic. *Theor. Comput. Sci.*, 50:1–102, 1987.

- [12] J. Gustavsson and J. Svenningsson. A usage analysis with bounded usage polymorphism and subtyping. In M. Mohnen and P. W. M. Koopman, editors, *Implementation of Functional Languages, 12th International Workshop, IFL 2000, Aachen, Germany, September 4-7, 2000, Selected Papers*, volume 2011 of *Lecture Notes in Computer Science*, pages 140–157. Springer, 2000.
- [13] J. C. Guzmán and P. Hudak. Single-threaded polymorphic lambda calculus. In *Proceedings of the Fifth Annual Symposium on Logic in Computer Science (LICS '90)*, Philadelphia, Pennsylvania, USA, June 4-7, 1990, pages 333–343, 1990.
- [14] J. Hage, S. Holdermans, and A. Middelkoop. A generic usage analysis with subeffect qualifiers. In R. Hinze and N. Ramsey, editors, *Proceedings of the 12th ACM SIGPLAN International Conference on Functional Programming, ICFP 2007, Freiburg, Germany, October 1-3, 2007*, pages 235–246. ACM, 2007.
- [15] S. Holdermans and J. Hage. Making “strictness” more relevant. *Higher-Order and Symbolic Computation*, 23(3):315–335, 2010.
- [16] K. Honda. Types for dyadic interaction. In *CONCUR*. Springer, 1993.
- [17] M. P. Jones. *Qualified Types: Theory and Practice*. Cambridge University Press, 1994.
- [18] M. P. Jones. Simplifying and improving qualified types. In *Proceedings of the seventh international conference on Functional programming languages and computer architecture, FPCA '95*, pages 160–169, La Jolla, California, USA, 1995. ACM.
- [19] M. P. Jones. First-class polymorphism with type inference. In P. Lee, F. Henglein, and N. D. Jones, editors, *Conference Record of POPL'97: The 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium, Paris, France, 15-17 January 1997*, pages 483–496. ACM Press, 1997.
- [20] M. P. Jones. Type classes with functional dependencies. In *Proceedings of the 9th European Symposium on Programming Languages and Systems, ESOP '00*, pages 230–244, Berlin, Germany, 2000. Springer-Verlag.
- [21] J. Launchbury and R. Paterson. Parametricity and unboxing with unpointed types. In H. R. Nielson, editor, *Programming Languages and Systems - ESOP'96, 6th European Symposium on Programming, Linköping, Sweden, April 22-24, 1996, Proceedings*, volume 1058 of *Lecture Notes in Computer Science*, pages 204–218. Springer, 1996.
- [22] O. Lee and K. Yi. Proofs about a folklore let-polymorphic type inference algorithm. *ACM Trans. Program. Lang. Syst.*, 20(4):707–723, 1998.
- [23] S. Lindley and J. G. Morris. A semantics for propositions as sessions. In J. Vitek, editor, *Programming Languages and Systems - 24th European Symposium on Programming, ESOP 2015, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2015, London, UK, April 11-18, 2015. Proceedings*, volume 9032 of *Lecture Notes in Computer Science*, pages 560–584. Springer, 2015.
- [24] K. Mazurak, J. Zhao, and S. Zdancewic. Lightweight linear types in System F°. In *TLDI*, 2010.
- [25] J. G. Morris. The best of both worlds: Linear functional programming without compromise (extended version). <http://homepages.inf.ed.ac.uk/jmorri14/>, 2016.
- [26] G. Morrisett, A. J. Ahmed, and M. Fluet. L³: A linear language with locations. In P. Urzyczyn, editor, *Typed Lambda Calculi and Applications, 7th International Conference, TLCA 2005, Nara, Japan, April 21-23, 2005, Proceedings*, volume 3461 of *Lecture Notes in Computer Science*, pages 293–307. Springer, 2005.
- [27] P. W. O’Hearn and D. J. Pym. The logic of bunched implications. *Bulletin of Symbolic Logic*, 5(2):215–244, 1999.
- [28] T. Petricek, D. A. Orchard, and A. Mycroft. Coeffects: a calculus of context-dependent computation. In J. Jeuring and M. M. T. Chakravarty, editors, *Proceedings of the 19th ACM SIGPLAN international conference on Functional programming, Gothenburg, Sweden, September 1-3, 2014*, pages 123–135. ACM, 2014.
- [29] R. Pucella and J. A. Tov. Haskell session types with (almost) no class. In *Proceedings of the 1st ACM SIGPLAN Symposium on Haskell, Haskell 2008, Victoria, BC, Canada, 25 September 2008*, pages 25–36. ACM, 2008.
- [30] J. C. Reynolds. Syntactic control of interference. In A. V. Aho, S. N. Zilles, and T. G. Szymanski, editors, *Conference Record of the Fifth Annual ACM Symposium on Principles of Programming Languages, Tucson, Arizona, USA, January 1978*, pages 39–46. ACM Press, 1978.
- [31] S. Smetsers, E. Barendsen, M. C. J. D. van Eekelen, and M. J. Plasmeijer. Guaranteeing safe destructive updates through a type system with uniqueness information for graphs. In H. J. Schneider and H. Ehrig, editors, *Graph Transformations in Computer Science, International Workshop, Dagstuhl Castle, Germany, January 1993, Proceedings*, volume 776 of *Lecture Notes in Computer Science*, pages 358–379. Springer, 1993.
- [32] F. Smith, D. Walker, and J. G. Morrisett. Alias types. In G. Smolka, editor, *Programming Languages and Systems, 9th European Symposium on Programming, ESOP 2000, Held as Part of the European Joint Conferences on the Theory and Practice of Software, ETAPS 2000, Berlin, Germany, March 25 - April 2, 2000, Proceedings*, volume 1782 of *Lecture Notes in Computer Science*, pages 366–381. Springer, 2000.
- [33] J. A. Tov and R. Pucella. Practical affine types. In T. Ball and M. Sagiv, editors, *Proceedings of the 38th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2011, Austin, TX, USA, January 26-28, 2011*, pages 447–458. ACM, 2011.
- [34] V. T. Vasconcelos, S. J. Gay, and A. Ravara. Type checking a multi-threaded functional language with session types. *Theor. Comput. Sci.*, 368(1-2):64–87, 2006.
- [35] P. Wadler. A taste of linear logic. In *Mathematical Foundations of Computer Science 1993, 18th International Symposium, MFCS'93, Gdansk, Poland, August 30 - September 3, 1993, Proceedings*, pages 185–210, 1993.
- [36] P. Wadler. Propositions as sessions. *J. Funct. Program.*, 24(2-3):384–418, 2014.
- [37] D. Walker. Substructural type systems. In B. C. Pierce, editor, *Advanced Topics in Types and Programming Languages*, chapter 1. The MIT Press, 2004.
- [38] D. Walker, K. Cray, and J. G. Morrisett. Typed memory management via static capabilities. *ACM Trans. Program. Lang. Syst.*, 22(4):701–771, 2000.